

BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections

Konstantin Vorontsov¹, Oleksandr Frei², Murat Apishev³, Peter Romov⁴, and Marina Dudarenko⁵

¹ Dorodnicyn Computing Centre of RAS, MIPT, HSE, Yandex, voron@forecsys.ru

² Schlumberger Information Solutions, oleksandr.frei@gmail.com

³ Lomonosov Moscow State University, great-mel@yandex.ru

⁴ Moscow Institute of Physics and Technology, Yandex, peter@romov.ru

⁵ Lomonosov Moscow State University, m.dudarenko@gmail.com

Abstract. Probabilistic topic modeling of text collections is a powerful tool for statistical text analysis. In this paper we announce the BigARTM open source project (<http://bigartm.org>) for regularized multimodal topic modeling of large collections. Several experiments on Wikipedia corpus show that BigARTM performs faster and gives better perplexity comparing to other popular packages, such as Vowpal Wabbit and Gensim. We also demonstrate several unique BigARTM features, such as additive combination of regularizers, topic sparsing and decorrelation, multimodal and multilanguage modeling, which are not available in the other software packages for topic modeling.

Keywords: probabilistic topic modeling, Probabilistic Latent Semantic Analysis, Latent Dirichlet Allocation, Additive Regularization of Topic Models, stochastic matrix factorization, EM-algorithm, BigARTM.

1 Introduction

Topic modeling is a rapidly developing branch of statistical text analysis [1]. Topic models reveal a hidden thematic structure of text collections and produce a compressed representation of documents in terms of their topics. Practical applications of topic models include many areas, such as information retrieval for long-text queries, classification, categorization, summarization and segmentation of texts. More ideas, models and applications are outlined in the survey [4].

From a statistical point of view, a probabilistic topic model defines each topic by a multinomial distribution over words, and then describes each document with a multinomial distribution over topics. From an optimizational point of view, topic modeling can be considered as a special case of approximate stochastic matrix factorization. To learn a factorized representation of a text collection is an ill-posed problem, which has an infinite set of solutions. A typical approach in this case is to apply regularization techniques, which impose problem-specific constraints and ultimately lead to a better solution.

Existing works on topic modeling describe hundreds of models, adapted to different situations. For practitioners, most of the models are too difficult to

quickly understand, compare, combine and embed into applications. This leads to a common practice of using only the basic models such as *Probabilistic Latent Semantic Analysis*, PLSA [6] and *Latent Dirichlet Allocation*, LDA [3]. Some of the difficulties are rooted in Bayesian learning framework, which is the dominating approach in topic modeling. Bayesian inference of topic models requires a laborious mathematical work, which prevents unification and flexible manipulation of various topic models. Until now, there was no freely available development tools to easily design, modify, select, and combine topic models.

In this paper we announce **the BigARTM open source project** for regularized multimodal topic modeling of large collections, <http://bigartm.org>. The theory behind BigARTM is based on a non-Bayesian multicriteria approach — *Additive Regularization of Topic Models*, ARTM [12]. In ARTM a topic model is learned by maximizing a weighted sum of the log-likelihood and additional regularization criteria. The optimization problem is solved by a general regularized expectation-maximization (EM) algorithm, which can be applied to an arbitrary combination of regularization criteria. Many known Bayesian topic models were revisited in terms of ARTM in [14,13]. Compared to the Bayesian approach, ARTM makes it easier to design, infer and combine topic models, thus reducing the barrier for entering into topic modeling research field.

BigARTM source code is released under the New BSD License, which permits free commercial and non-commercial usage. The core of the library is written in C++ and is exposed via two equally rich APIs for C++ and Python. The library is cross-platform and can be built for Linux, Windows and OS X in both 32 and 64 bit configuration. In our experiments on Wikipedia corpus BigARTM performs better than Vowpal Wabbit LDA and Gensim libraries in terms of perplexity and runtime. Comparing to the other libraries BigARTM offers several additional features, such as regularization and multi-modal topic modeling.

The rest of the paper is organized as follows. In section 2 we introduce a multimodal topic modeling for documents with additional discrete metadata. In section 3 we generalize the fast online algorithm [5] to additively regularized multimodal topic models. In section 4 we describe parallel architecture and implementation details of the BigARTM library. In section 5 we report results of our experiments on large datasets. In section 6 we discuss advantages, limitations and open problems of BigARTM.

2 Multimodal regularized topic model

Let D denote a finite set (collection) of texts and W^1 denote a finite set (vocabulary) of all terms from these texts. Each term can represent a single word or a key phrase. A document can contain not only words, but also terms of other modalities. Each modality is defined by a finite set (vocabulary) of terms W^m , $m = 1, \dots, M$. Examples of not-word modalities are: authors, class or category labels, date-time stamps, references to/from other documents, entities mentioned in texts, objects found in the images associated with the documents, users that read or downloaded documents, advertising banners, etc.

Assume that each term occurrence in each document refers to some latent topic from a finite set of topics T . Text collection is considered to be a sample of triples (w_i, d_i, t_i) , $i = 1, \dots, n$, drawn independently from a discrete distribution $p(w, d, t)$ over the finite space $W \times D \times T$, where $W = W^1 \sqcup \dots \sqcup W^M$ is disjoint a union of the vocabularies across all modalities. Terms w_i and documents d_i are observable variables, while topics t_i are latent variables.

Following the idea of Correspondence LDA [2] and Dependency LDA [10] we introduce a topic model for each modality:

$$p(w | d) = \sum_{t \in T} p(w | t) p(t | d) = \sum_{t \in T} \phi_{wt} \theta_{td}, \quad d \in D, w \in W^m, m = 1, \dots, M.$$

The parameters $\phi_{wt} = p(w | t)$ and $\theta_{td} = p(t | d)$ form matrices $\Phi^m = (\phi_{wt})_{W^m \times T}$ of *term probabilities for the topics*, and $\Theta = (\theta_{td})_{T \times D}$ of *topic probabilities for the documents*. The matrices Φ^m , if stacked vertically, form a $W \times T$ -matrix Φ . Matrices Φ^m and Θ are *stochastic*, that is, their vector-columns represent discrete distributions. The number of topics $|T|$ is expected to be much smaller than $|D|$ and $|W|$.

To learn parameters Φ^m, Θ from the multimodal text collection we maximize the log-likelihood for each m -th modality:

$$\mathcal{L}_m(\Phi^m, \Theta) = \sum_{d \in D} \sum_{w \in W^m} n_{dw} \ln p(w | d) \rightarrow \max_{\Phi^m, \Theta},$$

where n_{dw} is the number of occurrences of the term $w \in W^m$ in the document d . Note that topic distributions of documents Θ are common for all modalities. Following the ARTM approach, we add a regularization penalty term $R(\Phi, \Theta)$ and solve a constrained multicriteria optimization problem via scalarization:

$$\sum_{m=1}^M \tau_m \mathcal{L}_m(\Phi^m, \Theta) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad (1)$$

$$\sum_{w \in W^m} \phi_{wt} = 1, \phi_{wt} \geq 0; \quad \sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0. \quad (2)$$

The local maximum (Φ, Θ) of the problem (1), (2) satisfies the following system of equations with auxiliary variables $p_{tdw} = p(t | d, w)$:

$$p_{tdw} = \text{norm}_{t \in T}(\phi_{wt} \theta_{td}); \quad (3)$$

$$\phi_{wt} = \text{norm}_{w \in W^m} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad n_{wt} = \sum_{d \in D} n_{dw} p_{tdw}; \quad (4)$$

$$\theta_{td} = \text{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{w \in d} \tau_{m(w)} n_{dw} p_{tdw}; \quad (5)$$

where operator $\text{norm}_{t \in T} x_t = \frac{\max\{x_t, 0\}}{\sum_{s \in T} \max\{x_s, 0\}}$ transforms a vector $(x_t)_{t \in T}$ to a discrete distribution; $m(w)$ is the modality of the term w , so that $w \in W^{m(w)}$.

The system of equations (3)–(5) follows from Karush–Kuhn–Tucker conditions (see Appendix A for the proof). It can be solved by various numerical methods. Particularly, the simple-iteration method is equivalent to the EM algorithm, which is typically used in practice. For single modality ($M = 1$) it gives the regularized EM algorithm proposed in [12]. With no regularization ($R = 0$) it corresponds to PLSA [6].

Many Bayesian topic models can be considered as special cases of ARTM with different regularizers R , as shown in [14,13]. For example, LDA [3] corresponds to the entropy smoothing regularizer.

Due to the unified framework of additive regularization BigARTM can build topic models for various applications simply by choosing a suitable combination of regularizers from a build-in user extendable library.

3 Online topic modeling

Following the idea of Online LDA [5] we split the collection D into batches D_b , $b = 1, \dots, B$, and organize EM iterations so that each document vector θ_d is iterated until convergence at a constant matrix Φ , see Algorithm 1 and 2. Matrix Φ is updated rarely, after all documents from the batch are processed. For a large collection matrix Φ often stabilizes after small initial part of the collection. Therefore a single pass through the collection might be sufficient to learn a topic model. The second pass may be needed for the initial part of the collection.

Algorithm 1 does not specify how often to synchronize Φ matrix at steps 5–8. It can be done after every batch or less frequently (for instance if $\frac{\partial R}{\partial \phi_{wt}}$ takes long time to evaluate). This flexibility is especially important for concurrent implementation of the algorithm, where multiple batches are processed in parallel. In this case synchronization can be triggered when a fixed number of documents had been processed since the last synchronization.

The online reorganization of the EM iterations is not necessarily associated with Bayesian inference used in [5]. Different topic models, from PLSA to multi-modal and regularized models, can be learned by the above online EM algorithm.

4 BigARTM architecture

The main goal for BigARTM architecture is to ensure a constant memory usage regardless of the collection size. For this reason each D_b batch is stored on disk in a separate file, and only a limited number of batches is loaded into the main memory at any given time. The entire Θ matrix is also never stored in the memory. As a result, the memory usage stays constant regardless of the size of the collection.

Concurrency. An general rule of concurrency design is to express parallelism at the highest possible level. For this reason BigARTM implements a concurrent processing of the batches and keeps a single-threaded code for the `ProcessBatch(D_b, ϕ_{wt})` routine.

Algorithm 1: Online EM-algorithm for multimodal ARTM

Input: collection D_b , discounting factor $\rho \in (0, 1]$;
Output: matrix Φ ;

- 1 initialize ϕ_{wt} for all $w \in W$ and $t \in T$;
- 2 $n_{wt} := 0$, $\tilde{n}_{wt} := 0$;
- 3 **for all** batches D_b , $b = 1, \dots, B$
- 4 $\tilde{n}_{wt} := \tilde{n}_{wt} + \text{ProcessBatch}(D_b, \phi_{wt})$;
- 5 **if** (*synchronize*) **then**
- 6 $n_{wt} := \rho n_{wt} + \tilde{n}_{dw}$ for all $w \in W$ and $t \in T$;
- 7 $\phi_{wt} := \text{norm}_{w \in W^m} (n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}})$ for all $w \in W^m$, $m = 1, \dots, M$ and $t \in T$;
- 8 $\tilde{n}_{wt} := 0$ for all $w \in W$ and $t \in T$;

Algorithm 2: ProcessBatch(D_b, ϕ_{wt})

Input: batch D_b , matrix ϕ_{wt} ;
Output: matrix \tilde{n}_{wt} ;

- 1 $\tilde{n}_{wt} := 0$ for all $w \in W$ and $t \in T$;
- 2 **for all** $d \in D_b$
- 3 initialize $\theta_{td} := \frac{1}{|T|}$ for all $t \in T$;
- 4 **repeat**
- 5 $p_{tdw} := \text{norm}_{t \in T} (\phi_{wt} \theta_{td})$ for all $t \in T$;
- 6 $n_{td} := \sum_{w \in d} \tau_m(w) n_{dw} p_{tdw}$ for all $t \in T$;
- 7 $\theta_{td} := \text{norm}_{t \in T} (n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}})$ for all $t \in T$;
- 8 **until** θ_d converges;
- 9 increment \tilde{n}_{wt} by $n_{dw} p_{tdw}$ for all $w \in d$ and $t \in T$;

To split collection into batches and process them concurrently is a common approach, introduced in AD-LDA algorithm [8], and then further developed in PLDA [15] and PLDA+ [7] algorithms. These algorithms require all concurrent workers to become idle before an update of the Φ matrix. Such synchronization step adds a large overhead in the online algorithm where Φ matrix is updated multiple times on each iteration. An alternative architecture without the synchronization step is described in [11], however it mostly targets a distributed cluster environment. In our work we develop an efficient single-node architecture where all workers benefit from the shared memory space.

To run multiple ProcessBatch in parallel the inputs and outputs of this routine are stored in two separate in-memory queues, locked for push and pop operations with spin locks. This approach does not add any noticeable synchronization overhead because both queues only store smart pointers to the actual data objects, so push and pop operations does not involve copying or relocating big objects in the memory.

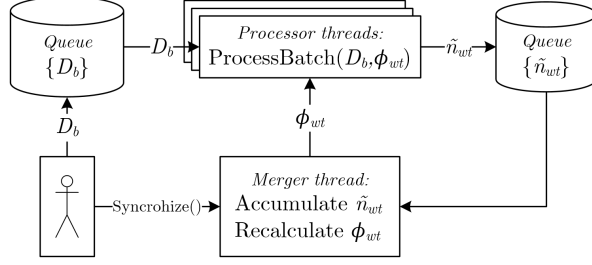


Fig. 1. Diagram of key BigARTM components

Smart pointers are also essential for lifecycle of the Φ matrix. This matrix is *read* by all processors threads, and can be *written* at any time by the merger thread. To update Φ without pausing all processor threads we keep two copies — an *active* Φ and a *background* Φ matrices. The active matrix is read-only, and is used by the processor threads. The background matrix is being built in a background by the merger thread at steps 6 and 7 of Algorithm 1, and once it is ready merger thread marks it as active. Before processing a new batch the processor thread gets the current active matrix from the merger thread. This object is passed via shared smart pointer to ensure that processor thread can keep ownership of its Φ matrix until the batch is fully processed. As a result, all processor threads keep running concurrently with the update of Φ matrix.

Note that all processor threads share the same Φ matrix, which means that memory usage stays at constant level regardless of how many cores are used for computation. Using memory for two copies of the Φ matrix in our opinion gives a reasonable usage balance between memory and CPU resources. An alternative solution with only one Φ matrix is also possible, but it would require a heavy usage of atomic CPU instructions. Such operations are very efficient, but still come at a considerable synchronization cost⁶, and using them for all reads and writes of the Φ matrix would cause a significant performance degradation for merger and processor threads. Besides, an arbitrary overlap between reads and writes of the Φ matrix eliminates any possibility of producing a deterministic result. The design with two copies of the Φ matrix gives much more control over this and in certain cases allows BigARTM to behave in a fully deterministic way.

The design with two Φ matrices only supports a single merger thread, and we believe it should handle all \tilde{n}_{wt} updates coming from many threads. This is a reasonable assumption because merging at step 6 takes only about $O(|W| \cdot |T|)$ operations to execute, while `ProcessBatch` takes $O(n|T|I)$ operations, where n is the number of non-zero entries in the batch, I is the average number of inner iterations in `ProcessBatch` routine. The ratio $n/|W|$ is typically from 100 to 1000 (based on datasets in UCI Bag-Of-Words repository), and I is 10...20, so the ratio safely exceeds the expected number of cores (up to 32 physical CPU cores in modern workstations, and even 60 cores of the Intel Xeon Phi co-processors).

⁶ <http://stackoverflow.com/questions/2538070/atomic-operation-cost>

Data layout. BigARTM uses dense single-precision matrices to represent Φ and Θ . Together with the Φ matrix we store a global dictionary of all terms $w \in W$. This dictionary is implemented as `std::unordered_map` that maps a string representation of $w \in W$ into its integer index in the Φ matrix. This dictionary can be extended automatically as more and more batches came through the system. To achieve this each batch D_b contains a local dictionary W_b , listing all terms that occur in the batch. The n_{dw} elements of the batch are stored as a sparse CSR matrix (Compressed Sparse Row format), where each row correspond to a document $d \in D_b$, and terms w run over a local batch dictionary W_b .

For performance reasons Φ matrix is stored in column-major order, and Θ in row-major order. This layout ensures that $\sum_t \phi_{wt} \theta_{td}$ sum runs on contiguous memory blocks. In both matrices all values smaller than 10^{-16} are always replaced with zero to avoid performance issues with denormalized numbers⁷.

Programming interface. All functionality of BigARTM is expressed in a set of extern C methods. To input and output complex data structures the API uses Google Protocol Buffers⁸. This approach makes it easy to integrate BigARTM into any research or production environment, as almost every modern language has an implementation of Google Protocol Buffers and a way of calling extern C code (ctypes module for Python, loadlibrary for Matlab, PInvoke for C#, etc).

On top of the extern C API BigARTM already has convenient wrappers in C++ and Python. We are also planning to implement a Java wrapper in the near future. In addition to the APIs the library also has a simple CLI interface.

BigARTM has built-in libraries of regularizers and quality measures that can be extended in current implementation only through project recompilation.

Basic tools. A careful selection of the programming tools is important for any software project. This is especially true for BigARTM as its code is written in C++, a language that by itself offers less functionality comparing to Python, .NET Framework or Java. To mitigate this we use various parts of the Boost C++ Libraries, Google Protocol Buffers for data serialization, ZeroMQ library for network communication, and several other libraries.

BigARTM uses CMake as a cross-platform build system, and it successfully builds on Windows, Linux and OS X in 32 and 64 bit configurations. Building the library require a recent C++ compiler with C++11 support (GNU GCC 4.6.3, clang 3.4 or Visual Studio 2012 or newer), and Boost Libraries 1.46.1 or newer. All the other third-parties are included in BigARTM repository.

We also use free online services to store source code (<https://github.com>), to host online documentation (<https://readthedocs.org>) and to run automated continuous integration builds (<http://travis-ci.org>).

⁷ http://en.wikipedia.org/wiki/Denormal_number#Performance_issues

⁸ <http://code.google.com/p/protobuf/>

5 Experiments

In this section we evaluate the runtime performance and the algorithmic quality of BigARTM against two popular software packages — Gensim [9] and Vowpal Wabbit⁹. We also demonstrate some of the unique BigARTM features, such as combining regularizers and multi-language topic modeling via multimodality, which are not available in the other software packages.

All three libraries (VW.LDA, Gensim and BigARTM) work out-of-core, e. g. they are designed to process data that is too large to fit into a computer’s main memory at one time. This allowed us to benchmark on a fairly large collection — 3.7 million articles from the English Wikipedia¹⁰. The conversion to bag-of-words was done with `gensim.make_wikicorpus` script¹¹, which excludes all non-article pages (such as category, file, template, user pages, etc), and also pages that contain less than 50 words. The dictionary is formed by all words that occur in at least 20 documents, but no more than in 10% documents in the collection. The resulting dictionary was capped at $|W| = 100\,000$ most frequent words.

Both Gensim and VW.LDA represents the resulting topic model as Dirichlet distribution over Φ and Θ matrices: $\theta_d \sim \text{Dir}(\gamma_d)$ and $\phi_t \sim \text{Dir}(\lambda_t)$. On contrary, BigARTM outputs a non-probabilistic matrices Φ and Θ . To compare the perplexity we take the maximum likelihood (ML) and maximum a posteriori (MAP) estimates from γ_d and λ_t distributions:

$$\begin{aligned}\theta_{td}^{\text{ML}} &\propto \gamma_{td} + \alpha, & \phi_{wt}^{\text{ML}} &\propto \lambda_{wt} + \beta; \\ \theta_{td}^{\text{MAP}} &\propto \max\{\gamma_{td} + \alpha - 1, 0\}, & \phi_{wt}^{\text{MAP}} &\propto \max\{\lambda_{wt} + \beta - 1, 0\}.\end{aligned}$$

The perplexity measure is defined as

$$\mathcal{P}(D, p) = \exp\left(-\frac{1}{n} \sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w | d)\right). \quad (6)$$

Comparison to existing software packages. The *Vowpal Wabbit* (VW) is a library of online algorithms that cover a wide range of machine learning problems. For topic modeling VW has the VW.LDA algorithm, based on the Online Variational Bayes LDA [5]. VW.LDA is neither multi-core nor distributed, but an effective single-threaded implementation in C++ made it one of the fastest tools for topic modeling.

The *Gensim* library specifically targets the area of topic modeling and matrix factorization. It has two LDA implementations — `LdaModel` and `LdaMulti-core`, both based on the same algorithm as VW.LDA (Online Variational Bayes LDA [5]). Gensim is entirely written in Python. Its high performance is achieved through the usage of NumPy library, built over low-level BLAS libraries (such as Intel MKL, ATLAS, or OpenBLAS). In `LdaModel` all batches are processed

⁹ https://github.com/JohnLangford/vowpal_wabbit/

¹⁰ <http://dumps.wikimedia.org/enwiki/20141208/>

¹¹ <https://github.com/piskvorky/gensim/tree/develop/gensim/scripts/>

Table 1. The comparison of BigARTM with VW.LDA and Gensim. *Train time* is the time for model training, *inference* is the time for calculation of θ_d of 100 000 held-out documents, *perplexity* is calculated according to (6) on held-out documents.

library	procs	train time	inference	perplexity	
				mean	map
BigARTM +smoothing	1	62 min	127 sec	4000	
Gensim LDA	1	369 min	395 sec	4161	4213
Vowpal Wabbit LDA	1	73 min	120 sec	4108	4061
BigARTM +smoothing	4	13 min	33 sec	4061	
Gensim LDA-Multicore	4	60 min	222 sec	4111	4055
BigARTM +smoothing	8	8 min	24 sec	4304	
Gensim LDA-Multicore	8	57 min	224 sec	4455	4379

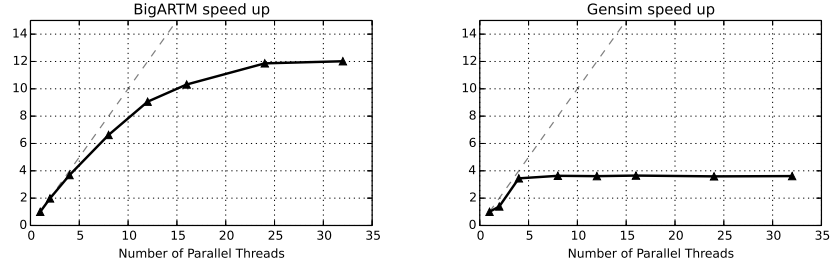


Fig. 2. Speed up for BigARTM (left chart) and Gensim (right chart)

sequentially, and the concurrency happens entirely within NumPy. In LdaMulticore the workflow is similar to BigARTM — several batches are processed concurrently, and there is a single aggregation thread that asynchronously merges the results.

Each run in our experiment performs one pass over the Wikipedia corpus and produces a model with $|T| = 100$ topics. The runtime is reported for an Intel-based CPU with 16 physical cores with hyper-threading. The collection was split into batches with 10000 documents each (`chunksize` in Gensim, `minibatch` in VW.LDA). The update rule in online algorithm used $\rho = (b + \tau_0)^{-0.5}$, where b is the number of batches processed so far, and τ_0 is an a constant offset parameter introduced in [5], in our experiment $\tau_0 = 64$. Updates were performed after each batch in non-parallel runs, and after P batches when running in P threads. LDA priors were fixed as $\alpha = 0.1$, $\beta = 0.1$, so that $\theta_d \sim \text{Dir}(\alpha)$, $\phi_t \sim \text{Dir}(\beta)$.

Table 1 compares the performance of VW.LDA, Gensim, and BigARTM.

Fig. 2 compares BigARTM and Gensim speedup depending on the number of CPU threads for Amazon AWS c3.8xlarge with 32 cores, Gensim 0.10.3 under Python 2.7, and the parameter of LDA-Multicore `workers = nProcessors - 1`.

Experiments with combination of regularizers. BigARTM has a built-in library of regularizers, which can be used in any combination. In the following experiment we combine three regularizers: sparsing of ϕ_t distributions, sparsing of θ_d distri-

Table 2. Comparison of LDA and ARTM models. Quality measures: \mathcal{P}_{10k} , \mathcal{P}_{100k} — hold-out perplexity on 10K and 100K documents sets, \mathcal{S}_Φ , \mathcal{S}_Θ — sparsity of Φ and Θ matrices (in %), \mathcal{K}_s , \mathcal{K}_p , \mathcal{K}_c — average topic kernel size, purity and contrast respectively.

Model	\mathcal{P}_{10k}	\mathcal{P}_{100k}	\mathcal{S}_Φ	\mathcal{S}_Θ	\mathcal{K}_s	\mathcal{K}_p	\mathcal{K}_c
LDA	3499	3827	0.0	0.0	931	0.535	0.516
ARTM	3592	3944	96.3	80.5	1135	0.810	0.732

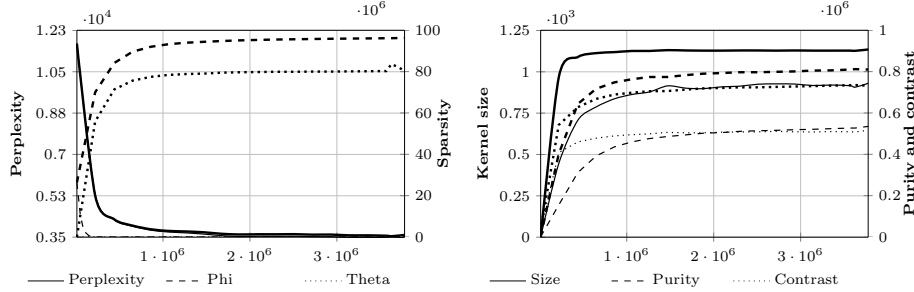


Fig. 3. Comparison of LDA (thin) and ARTM (bold) models. The number of processed documents is shown along the X axis.

butions, and pairwise decorrelation of ϕ_t distributions. This combination helps to improve several quality measures without significant loss of perplexity, according to experiments on the offline implementation of ARTM [14]. The goal of our experiment is to show that this remains true for the online implementation in BigARTM. We use the following built-in quality measures: the hold-out perplexity, the sparsity of Φ and Θ matrices, and the characteristics of topic lexical kernels (size, purity, and contrast) averaged across all topics.

Table 2 compares the results of additive combination of regularizers (ARTM) and the usual LDA model. Figure 3 presents quality measures as functions of the number of processed documents. The left chart shows perplexity and sparsity of Φ , Θ matrices, and the right chart shows average lexical kernel measures.

Experiments on multi-language Wikipedia. To show how BigARTM works with multimodal datasets we prepared a text corpus containing all English and Russian Wikipedia articles with mutual interlanguage links. We represent each linked pair of articles as a single multi-language document with two modalities, one modality for each language. That is how our multi-language collection acts as a multimodal document collection.

The dump of Russian articles¹² had been processed following the same technique as we previously used in experiments on English Wikipedia. Russian words were lemmatized with Yandex MyStem 3.0¹³. To further reduce the dictionary we only keep words that appear in no less than 20 documents, but no more

¹² <http://dumps.wikimedia.org/ruwiki/20141203/>

¹³ <https://tech.yandex.ru/mystem/>

Table 3. Top 10 words with $p(w|t)$ probabilities (in %) from two-language topic model, based on Russian and English Wikipedia articles with mutual interlanguage links.

Topic 68				Topic 79			
research	4.56	институт	6.03	goals	4.48	матч	6.02
technology	3.14	университет	3.35	league	3.99	игрок	5.56
engineering	2.63	программа	3.17	club	3.76	сборная	4.51
institute	2.37	учебный	2.75	season	3.49	фк	3.25
science	1.97	технический	2.70	scored	2.72	против	3.20
program	1.60	технология	2.30	cup	2.57	клуб	3.14
education	1.44	научный	1.76	goal	2.48	футболист	2.67
campus	1.43	исследование	1.67	apps	1.74	гол	2.65
management	1.38	наука	1.64	debut	1.69	забивать	2.53
programs	1.36	образование	1.47	match	1.67	команда	2.14
Topic 88				Topic 251			
opera	7.36	опера	7.82	windows	8.00	windows	6.05
conductor	1.69	оперный	3.13	microsoft	4.03	microsoft	3.76
orchestra	1.14	дирижер	2.82	server	2.93	версия	1.86
wagner	0.97	певец	1.65	software	1.38	приложение	1.86
soprano	0.78	певица	1.51	user	1.03	сервер	1.63
performance	0.78	театр	1.14	security	0.92	server	1.54
mozart	0.74	партия	1.05	mittchell	0.82	программный	1.08
sang	0.70	сопрано	0.97	oracle	0.82	пользователь	1.04
singing	0.69	вагнер	0.90	enterprise	0.78	обеспечение	1.02
operas	0.68	оркестр	0.82	users	0.78	система	0.96

than in 10% of documents in the collection. The resulting collection contains 216175 pairs of Russian–English articles, with combined dictionary of 196749 words (43% Russian, 57% English words).

We build multi-language model with 400 topics. They cover a wide range of themes such as science, architecture, history, culture, technologies, army, different countries. All 400 topics were reviewed by an independent assessor, and he successfully interpreted all except four topics.

Table 3 shows top 10 words for four randomly selected topics. Top words in these topics are clearly consistent between Russian and English languages. The Russian part of last topic contains some English words such as “Windows” or “Server” because it is common to use them in Russian texts without translation.

6 Conclusions

BigARTM in an open source project for parallel online topic modeling of large text collections. It provides a high flexibility for various applications due to multimodality and additive combinations of regularizers. BigARTM architecture has a rich potential. Current components can be reused in a distributed solution that runs on cluster. Further improvement of single-node can be achieved by offloading batch processing into GPU.

Acknowledgements. The work was supported by the Russian Foundation for Basic Research grants 14-07-00847, 14-07-00908, 14-07-31176 and by Skolkovo Institute of Science and Technology (project 081-R).

References

1. Blei, D.M.: Probabilistic topic models. *Communications of the ACM* 55(4), 77–84 (2012)
2. Blei, D.M., Jordan, M.I.: Modeling annotated data. In: *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval*. pp. 127–134. ACM, New York, NY, USA (2003)
3. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet allocation. *Journal of Machine Learning Research* 3, 993–1022 (2003)
4. Daud, A., Li, J., Zhou, L., Muhammad, F.: Knowledge discovery through directed probabilistic topic models: a survey. *Frontiers of Computer Science in China* 4(2), 280–301 (2010)
5. Hoffman, M.D., Blei, D.M., Bach, F.R.: Online learning for latent dirichlet allocation. In: *NIPS*. pp. 856–864. Curran Associates, Inc. (2010)
6. Hofmann, T.: Probabilistic latent semantic indexing. In: *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*. pp. 50–57. ACM, New York, NY, USA (1999)
7. Liu, Z., Zhang, Y., Chang, E.Y., Sun, M.: PLDA+: Parallel latent dirichlet allocation with data placement and pipeline processing. *ACM Transactions on Intelligent Systems and Technology (TIST)* 2(3), 26 (2011)
8. Newman, D., Asuncion, A., Smyth, P., Welling, M.: Distributed algorithms for topic models. *The Journal of Machine Learning Research* 10, 1801–1828 (2009)
9. Řehůřek, R., Sojka, P.: Software Framework for Topic Modelling with Large Corpora. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. pp. 45–50. ELRA, Valletta, Malta (May 2010).
10. Rubin, T.N., Chambers, A., Smyth, P., Steyvers, M.: Statistical topic models for multi-label document classification. *Machine Learning* 88(1-2), 157–208 (2012)
11. Smola, A., Narayanamurthy, S.: An architecture for parallel topic models. *Proceedings of the VLDB Endowment* 3(1-2), 703–710 (2010)
12. Vorontsov, K.V.: Additive regularization for topic models of text collections. *Doklady Mathematics* 89(3), 301–304 (2014)
13. Vorontsov, K.V., Potapenko, A.A.: Additive regularization of topic models. *Machine Learning, Special Issue on Data Analysis and Intelligent Optimization* (2014)
14. Vorontsov, K.V., Potapenko, A.A.: Tutorial on probabilistic topic modeling: Additive regularization for stochastic matrix factorization. In: *AIST’2014, Analysis of Images, Social networks and Texts*. vol. 436, pp. 29–46. Springer International Publishing Switzerland, Communications in Computer and Information Science (CCIS) (2014)
15. Wang, Y., Bai, H., Stanton, M., Chen, W.Y., Chang, E.Y.: PLDA: Parallel latent dirichlet allocation for large-scale applications. In: *Algorithmic Aspects in Information and Management*, pp. 301–314. Springer (2009)

Appendix A

Consider the system of equations (3)–(5).

Topic t is called *regular* for the modality m if $n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}^m} > 0$ for at least one term $w \in W^m$. If the reverse inequality holds for all $w \in W^m$ then topic t is called *irregular*.

Document d is called *regular* if $n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} > 0$ for at least one topic $t \in T$. If the reverse inequality holds for all $t \in T$ then document d is called *irregular*.

Theorem 1. *If the function $R(\Phi, \Theta)$ is continuously differentiable and (Φ, Θ) is the local maximum of the problem (1), (2) then for any regular topic-modality pair (t, m) and any regular document d the system of equations (3)–(5) holds.*

Note 1. If a topic t is irregular then the t -th vector-column in matrix Φ^m equals zero and can not represent a discrete distribution. This means that topic t for the modality m must be excluded from the model. This mechanism is useful for irrelevant topics elimination and determining the number of topics.

Note 2. If a documents d is irregular then the d -th vector-column in matrix Θ equals zero and can not represent a discrete distribution. This means that document d must be excluded from the model. For example, a document may be too short or irrelevant to the given collection.

Proof. For the local minimum Φ^m, Θ of the problem (1), (2) the Karush–Kuhn–Tucker (KKT) conditions can be written as follows:

$$\begin{aligned} \sum_d n_{dw} \frac{\theta_{td}}{p(w|d)} + \frac{\partial R}{\partial \phi_{wt}} &= \lambda_t - \lambda_{wt}; \quad \lambda_{wt} \geq 0; \quad \lambda_{wt} \phi_{wt} = 0; \\ \sum_m \tau_m \sum_{w \in W^m} n_{dw} \frac{\phi_{wt}}{p(w|d)} + \frac{\partial R}{\partial \theta_{td}} &= \mu_d - \mu_{td}; \quad \mu_{td} \geq 0; \quad \mu_{td} \theta_{td} = 0; \end{aligned}$$

where $\lambda_t, \lambda_{wt}, \mu_d, \mu_{td}$ are KKT multipliers for normalization and nonnegativity constrains.

Let us multiply both sides of the first equation by ϕ_{wt} , both sides of the second equation by θ_{td} , and reveal the auxiliary variable p_{tdw} from (3) in the left-hand side of both equations. Then we sum the right-hand side of the first equation over d , the right-hand side of the second equation over t :

$$\begin{aligned} \phi_{wt} \lambda_t &= \sum_d n_{dw} \frac{\phi_{wt} \theta_{td}}{p(w|d)} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} = n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}; \\ \theta_{td} \mu_d &= \sum_m \tau_m \sum_{w \in W^m} n_{dw} \frac{\phi_{wt} \theta_{td}}{p(w|d)} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} = n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}. \end{aligned}$$

An assumption that $\lambda_t \leq 0$ contradicts the regularity condition for the (t, m) pair. Then $\lambda_t > 0$. Either $\phi_{wt} = 0$ or both sides of the first equation are positive. Combining these two cases in one formula, we write:

$$\phi_{wt} \lambda_t = \max \left\{ n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}, 0 \right\}. \quad (7)$$

Analogously, an assumption that $\mu_d \leq 0$ contradicts the regularity condition for the document d . Then $\mu_d > 0$. Either $\theta_{td} = 0$ or both sides of the second equation are positive, consequently,

$$\theta_{td}\mu_d = \max\left\{n_{td} + \theta_{td}\frac{\partial R}{\partial \theta_{td}}, 0\right\}. \quad (8)$$

Let us sum both sides of the first equation over all $w \in W^m$, then both sides of the second equation over all $t \in T$:

$$\lambda_t = \sum_{w \in W^m} \max\left\{n_{wt} + \phi_{wt}\frac{\partial R}{\partial \phi_{wt}}, 0\right\}; \quad (9)$$

$$\mu_d = \sum_{t \in T} \max\left\{n_{td} + \theta_{td}\frac{\partial R}{\partial \theta_{td}}, 0\right\}. \quad (10)$$

Finally, we obtain (4) by expressing ϕ_{wt} from (7) and (9).

Analogously, we obtain (5) by expressing θ_{td} from (8) and (10). \square